

Performance analysis of the Kahan-enhanced scalar product on current multi- and manycore processors

J. Hofmann^{1,*}, D. Fey¹, M. Riedmann², J. Eitzinger³, G. Hager³ and G. Wellein³

¹ Chair for Computer Architecture, University of Erlangen-Nuremberg, Erlangen, Germany

² AREVA GmbH, Erlangen, Germany

³ Erlangen Regional Computing Center (RRZE), University of Erlangen-Nuremberg, Erlangen, Germany

SUMMARY

We investigate the performance characteristics of a numerically enhanced scalar product (dot) kernel loop that uses the Kahan algorithm to compensate for numerical errors, and describe efficient SIMD-vectorized implementations on recent multi- and manycore processors. Using low-level instruction analysis and the execution-cache-memory (ECM) performance model we pinpoint the relevant performance bottlenecks for single-core and thread-parallel execution, and predict performance and saturation behavior. We show that the Kahan-enhanced scalar product comes at almost no additional cost compared to the naive (non-Kahan) scalar product if appropriate low-level optimizations, notably SIMD vectorization and unrolling, are applied. The ECM model is extended appropriately to accommodate not only modern Intel multicore chips but also the Intel Xeon Phi “Knights Corner” coprocessor and an IBM POWER8 CPU. This allows us to discuss the impact of processor features on the performance across four modern architectures that are relevant for high performance computing.

Copyright © 2016 John Wiley & Sons, Ltd.

Received ...

KEY WORDS: ECM Performance Model; Kahan; Scalar Product; Xeon; Knights Corner; POWER8

1. INTRODUCTION AND RELATED WORK

Accumulating finite-precision floating-point numbers in a scalar variable is a common operation in computational science and engineering. The consequences in terms of accuracy are inherent to the number representation and have been well known and studied for a long time [1]. There are a number of summation algorithms that enhance accuracy while maintaining an acceptable throughput [2, 3], of which Kahan [4] is probably the most popular one. However, the topic is still subject to active research [5, 6, 7, 8]. A straightforward solution to the inherent accuracy problems is arbitrary-precision floating point arithmetic, which comes at a significant performance penalty. Naive summation and arbitrary precision arithmetic are at opposite ends of a broad spectrum of options, and balancing performance vs. accuracy is a key concern when selecting a specific solution.

Naive summation, which simply adds each successive number in sequence to an accumulator, requires appropriate unrolling for single instruction multiple data (SIMD) vectorization and pipelining. The necessary code transformations are performed automatically by modern compilers,

[†]E-mail: johannes.hofmann@fau.de

*Correspondence to: Johannes Hofmann, Lehrstuhl für Rechnerarchitektur (Informatik 3), Martensstr. 3, 91058 Erlangen, Germany

which results in optimal in-core performance. Such a code quickly saturates the memory bandwidth of modern multi-core CPUs when the data is in memory.

This paper investigates implementations of the scalar product, a kernel which is relevant in many numerical algorithms. Starting from an optimal naive implementation it considers scalar and SIMD-vectorized versions of the Kahan algorithm using various SIMD instruction set extensions on a range of multi- and manycore processors from Intel and IBM. Using an analytic performance model we point out the conditions under which Kahan comes for free, and we predict the single core performance in all memory hierarchy levels as well as the scaling behavior across the cores of a chip. The present work is an extended version of [9], where we carried out the analysis for a range of older Intel Xeon processors. Apart from new architectures we present a refined version of the ECM performance model and add an additional optimization for the Intel Haswell-EP CPU.

This paper is organized as follows. In Sect. 3 we give an overview of the hardware used for analysis and benchmarking. Section 2 introduces the execution-cache-memory (ECM) performance model, which is used in Sect. 4 to describe different variants of the naive and the Kahan scalar product. Section 5 gives performance results and validates the models. Section 6 provides a conclusion and some comments on the possible extension of our work.

2. THE ECM PERFORMANCE MODEL

The execution-cache-memory (ECM) model [10, 11, 12, 9] is an analytic performance model that uses hardware architecture specifications and few measurements as input. It estimates the number of CPU cycles required to execute a number of iterations of a loop on a single core of a multi- or many-core chip. The prediction comprises contributions from the in-core execution time T_{core} , i.e., the time spent executing instructions in the core under the assumption that all data resides in the L1 cache, and the transfer time T_{data} , i.e., the time spent transferring data from its location in the cache/memory hierarchy to the L1 cache. As data transfers in the cache and memory hierarchy occur at cache line (CL) granularity we choose the number of loop iterations n_{it} to correspond to one cache line's "worth of work." On Intel architectures, where CLs are 64 B long, we use $n_{\text{it}} = 16$ for the single precision (SP) dot product because sixteen SP floating-point numbers (4 B each) fit into one CL. CLs on the IBM POWER8 architectures are 128 B, which leads to $n_{\text{it}} = 32$ for the SP dot product.

Superscalar core designs house multiple execution units, for loading and storing data, multiplication, division, addition, etc. The in-core execution time T_{core} is determined by the unit that takes the longest to execute the instructions allocated to it. Other constraints for the in-core execution time of a single core may apply, e.g., the four micro-op per cycle retirement limit for Intel's Xeon cores and the eight instruction per cycle retirement limit for IBM's POWER8 core. The model differentiates between core cycles depending on whether data transfers in the cache hierarchy can overlap with in-core execution time. For instance, on Intel Xeons, core cycles in which data is moved between the L1 cache and registers, e.g., cycles in which load and/or store instructions are retired, prohibit the simultaneous transfers of data between the L1 and L2 cache; these "non-overlapping" cycles contribute to T_{nOL} . Cycles in which other instructions, such as arithmetic instructions, retire are considered "overlapping" cycles and contribute to T_{OL} . The in-core runtime is the maximum of both: $T_{\text{core}} = \max(T_{\text{OL}}, T_{\text{nOL}})$. Note that the non-overlapping quality of L1-register transfers is specific to Intel CPUs. We will see later that the IBM POWER8 does not have non-overlapping instructions.

For modeling the data transfers, latency effects are initially neglected, so transfer times are exclusively a function of bandwidth. Cache bandwidths are typically well documented and can be found in vendor data sheets. Depending on how many CLs have to be transferred, the contribution of each level in the memory hierarchy ($T_{\text{L1L2}}, \dots, T_{\text{L3Mem}}$) can be determined. Special care has to be taken when dealing with main memory bandwidth, because peak memory bandwidth specified in the data sheet and sustained memory bandwidth b_s can differ greatly. In addition, in practice the sustained bandwidth may also depend on the number of distinct load and store streams. It is therefore recommended to empirically determine b_s using a kernel that resembles the memory access pattern

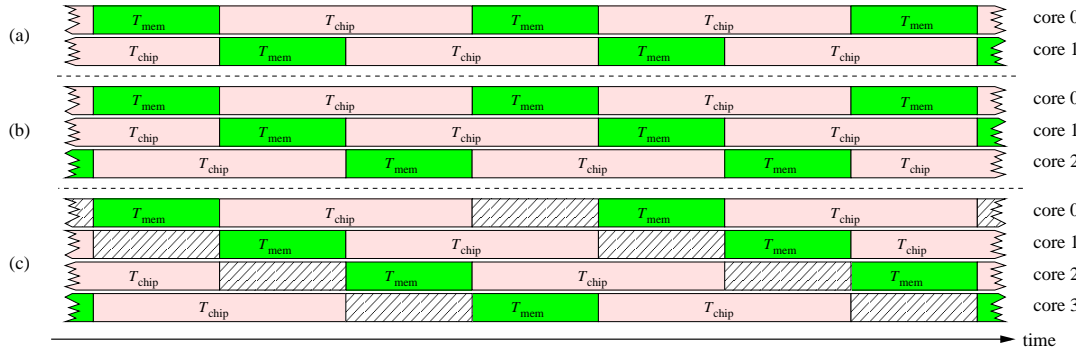


Figure 1. Multicore scaling in the ECM model. T_{mem} is the time spent for data transfers over the memory bottleneck, while T_{chip} comprises all non-bottlenecked (i.e., core-local) contributions. The saturation point is at three cores in this case since $(T_{\text{chip}} + T_{\text{mem}})/T_{\text{mem}} = 3$. Hatched boxes denote stalls, which emerge from using more cores than needed for bandwidth saturation.

of the benchmark to be modeled. Once b_s has been obtained, the time to transfer one CL between the cache hierarchy and main memory can be derived from the CPU frequency f as $64B \cdot f/b_s$ cycles.

In a second step, an empirically determined latency penalty T_p is applied to off-core transfer times. This departure from the bandwidth-only model has been mandated by the inability of some architectures to hide the memory access latency. On regular Xeon processors, this penalty is added for each level in the memory hierarchy that has to make use of the Uncore interconnect (i.e., the L3 cache, as data is pseudo-randomly distributed between all last-level cache segments and memory, because the memory controller is attached to the ring bus as well). On Knights Corner there exists no shared cache when each thread is working on its own data; each core is using data from its local L2 cache so the latency penalty is only added when the core-ring-interconnect is used to get data from main memory. Instruction times as well as data transfer times, e.g., T_{L1L2} for the time required to transfer data between L1 and L2 caches, are summarized in a shorthand notation: $\{T_{OL} \parallel T_{nOL} \mid T_{L1L2} \mid T_{L2L3} + T_p \mid T_{L3Mem} + T_p\}$.

To arrive at a prediction, in-core execution and data transfer times must be combined appropriately. The runtime is given by either T_{OL} or the sum of non-overlapping core cycles T_{nOL} plus contributions of data transfers T_{data} , whichever takes longer. T_{data} comprises all necessary data transfers in the memory hierarchy, plus latency penalties if applicable. Again we have to distinguish between overlapping and non-overlapping behavior; in case of Intel Xeon, any data transfer during a specific cycle in the inclusive cache hierarchy prevents all other transfers (including those in T_{OL}) in that cycle. T_{data} is thus the sum of all cycles required to transfer the data to L1 and back. E.g., for data coming from the L3 cache we have $T_{\text{data}} = T_{L1L2} + T_{L2L3} + T_p$. The prediction is thus $T_{\text{ECM}} = \max(T_{OL}, T_{nOL} + T_{\text{data}})$. Note that for other architectures with different overlapping properties and/or exclusive cache hierarchies this formula may look very different.

In order to summarize the predictions for data coming from different levels in the hierarchy we use a shorthand notation: $\{T_{\text{ECM}}^{\text{core}} \mid T_{\text{ECM}}^{\text{L2}} \mid T_{\text{ECM}}^{\text{L3}} \mid T_{\text{ECM}}^{\text{Mem}}\}$. Converting from time (cycles) to performance is done by dividing the work W (e.g., floating-point operations, updates, or any other relevant work metric) by the runtime: $P_{\text{ECM}} = W/T_{\text{ECM}}$.

The model assumes that single-core performance scales linearly with the cores until a shared bottleneck is saturated. On most modern processors the only shared bottleneck is main memory bandwidth. As shown in Fig. 1, the ratio of the overall single-core execution time and the contribution of the bottleneck determines the maximum speedup: as long as the number of cores is smaller than this ratio, the memory bus is not saturated. In terms of the ECM model, the maximum speedup is $\sigma_S = T_{\text{ECM}}^{\text{Mem}}/T_{L3\text{Mem}}$. Performance at the saturation point is then $P_{\text{ECM}}^S = f \cdot \sigma_S \cdot W_{\text{CL}}/T_{\text{ECM}}^{\text{Mem}} = f \cdot W_{\text{CL}}/T_{L3\text{Mem}}$, where W_{CL} is the work per CL and f is the processor clock frequency. This is just another formulation of the bandwidth-bound part of the Roofline model [13]. The core count necessary to saturate the memory bandwidth is $n_S = \lceil T_{\text{ECM}}^{\text{Mem}}/T_{L3\text{Mem}} \rceil$. If $n_S \geq n_{\text{chip}}$, i.e., if the required number of cores for saturation exceeds the available number, the code is scalable.

Microarchitecture	Haswell-EP	Broadwell-EP	Knights Corner	POWER8
Shorthand	HSW	BDW	KNC	PWR8
Chip model	E5-2695 v3	unknown	5110P	S822LC
Release date	Q3 2014	pre-release	Q4 2012	Q2 2014
Nominal CPU clock	2.3 GHz	2.1 GHz	1.05 GHz	2.926 GHz
Cores/threads	14/28	22/44	60/240	10/80
Max. SIMD width	32 B	32 B	64 B	16 B
# of SIMD registers	16	16	32	64
Instruction throughput per cycle				
LOAD/STORE	2 / 1	2 / 1	1 / 1	2 / 2
ADD/MUL/FMA	1 / 2 / 2	1 / 2 / 2	1 / 1 / 1	2 / 2 / 2
Core-private caches	32 kB L1	32 kB L1	32 kB L1	64 kB L1
	256 kB L2	256 kB L2	512 kB L2	512 kB L2
	—	—	—	8 MB L3
Shared caches	35 MB L3	55 MB L3	—	64 MB L4
L2-L1 bandwidth	64 B/cy	64 B/cy	32 B/cy	64 B/cy
L3-L2 bandwidth	32 B/cy	32 B/cy	—	32 B/cy
MEM-L3 bandwidth	~14 B/cy	~15 B/cy	160 B/cy	—
Centaur-L2 bandwidth	—	—	—	~19 B/cy
Main memory	4×DDR4-2166	4×DDR4-2166	16×GDDR5-5000	4×Centaur
Theor. load BW	69.3 GB/s	69.3 GB/s	320 GB/s	76.8 GB/s
Meas. load BW	2×32.0 GB/s (92%)	2×32.3 GB/s (93%)	175 GB/s (55%)	73.6 GB/s (96%)

Table I. Test machine specifications and micro-architectural features (one socket). The cache line size is 64 bytes for all Intel architectures and 128 bytes for IBM POWER8.

3. EXPERIMENTAL TESTBED

Table I gives an overview of the relevant architectural details of the systems used in this paper. The regular Xeon machines (Haswell-EP and Broadwell-EP) and the POWER8 machine are standard two-socket systems. The Xeon Phi coprocessor (Knights Corner) is a PCIe card hosted in a standard two-socket Ivy Bridge-EP system.

Note that BDW corresponds to a “tick” in Intel’s design model, i.e., it is a shrink in the manufacturing process technology from 22 nm to 14 nm with only minor architectural improvements compared to HSW. All results for Broadwell-EP are preliminary since we only had access to a pre-release version of the chip.

All SIMD instructions set extensions for the covered microarchitectures support fused multiply-add (FMA) instructions. The vector scalar extension (VSX) on IBM’s PWR8 have a SIMD width of 16 B. The AVX2 vector extensions supported by HSW and BDW have a SIMD width of 32 B and KNC’s initial many core instructions (IMCI) allow for 64-B SIMD. All Intel processors employ a fully inclusive cache architecture whereas PWR8 uses an exclusive victim cache architecture for the last level cache. This results in different data paths inside the caches. On PWR8 data is loaded from memory directly into the L2 caches, and only cachelines which get evicted from L2 will be copied back to the L3 cache.

The sustained memory bandwidth for all architectures was determined using a naive dot product benchmark. To obtain good results on the Xeon Phi, we followed the optimization instructions for the STREAM benchmark as described by Intel [14]; in particular, we set the prefetching distance 64 CLs ahead for the L2 cache, 8 CLs for the L1 cache, and used one thread per core to avoid

```

(a)
float sum = 0.0;
for (int i=0; i<N; i++) {
    sum = sum + a[i] * b[i]
}

(b)
float sum = 0.0;
float c = 0.0;
for (int i=0; i<N; ++i) {
    float prod = a[i]*b[i];
    float y = prod-c;
    float t = sum+y;
    c = (t-sum)-y;
    sum = t;
}

```

Figure 2. (a) Naive scalar product code in single precision. (b) Kahan-compensated scalar product code.

congestion on the ring bus. The IBM PWR8 memory bandwidth requires further explanation. PWR8 uses a custom high frequency channel interface between the processor chip and a memory buffer chip (Centaur) [15]. Each Centaur chip connects to four DRAM channels. PWR8 supports up to eight memory channels per chip operating at 9.6 GHz with a bus width of 2 B (read) plus 1 B (write). Our test system is an IBM S822LC and supports only four Centaur chips. The four memory channels can provide up to 115.2 GB/s read/write or 76.8 GB/s read-only bandwidth per chip. Note that this is significantly less than what the 16 attached DRAM channels (DDR3-1333) could provide (170.6 GB/s). A fully equipped high-end PWR8 system hence has twice the memory bandwidth per chip.

Unless noted otherwise, KNC was used in 2-SMT and PWR8 in 8-SMT mode, i.e., two respectively eight threads were run on each physical core. On HSW and BDW a single thread was run on each physical core, and Uncore frequency scaling was deactivated. Furthermore, the “cluster on die” (COD) mode was active for HSW and BDW. In CoD mode, the chip is logically split into two ccNUMA domains of equal size. Last-level cache and memory requests are limited to the domain a core is assigned to, reducing latency and collisions in the Uncore interconnect. The two memory domains per chip are visible in the load-only bandwidth row of Table I; e.g., the sustained load-only bandwidth for HSW is 32.0 GB/s per memory domain and 64.0 GB/s per chip. For details on the CoD mode see [16].

4. OPTIMAL IMPLEMENTATIONS AND PERFORMANCE MODELS FOR DOT

We only discuss variants for dot in SP here. The model prediction in terms of cycles per CL does not change for the SIMD variants of Kahan when going from SP to double precision (DP), but one CL update represents twice as much useful work (scalar iterations) in the SP case. To eliminate variations introduced by compiler-generated code we implemented all kernels directly in assembly language and use the *likwid-bench* microbenchmarking framework [17] to perform measurements.

4.1. Naive scalar product

An optimal implementation of the naive scalar product in single precision serves as the baseline (see Fig. 2a). All versions of the Kahan-enhanced scalar product described in Section 4.2 will be compared to this baseline.

Sufficient unrolling must be applied to hide the ADD pipeline latency for the recursive update on the accumulation register and to apply SIMD vectorization. Both optimizations introduce partial sums and are therefore not compatible with the C standard as the order of non-associative operations is changed. With higher optimization levels (`-O3`) the current Intel C compiler (version 15.0.2) and IBM XL C compiler (version 13.1.3) both generate optimal code. Note that partial sums usually improve the accuracy of the result [8].

4.1.1. Intel Haswell-EP and Broadwell-EP On HSW and BDW the kernel is limited by the throughput of the LOAD units (see Table I). Two AVX loads per vector (a and b) are required to cover one unit of work (16 scalar loop iterations), leading to a total of four AVX load instructions; with two LOAD units, the core can execute two LOAD instructions per cycle, resulting in $T_{nOL} = 2$ cy. To process the data, two FMA instructions have to be executed; with two FMA units, the core can execute both instructions in a single cycle, resulting in an overlapping part of $T_{OL} = 1$ cy.

If the data is in the L2 cache, two CLs (one each for a and b) have to be transferred to the L1 cache; at the advertised bandwidth of 64 B/cy this results in $T_{L1L2} = 2$ cy. With data in L3 it takes $T_{L2L3} = 4$ cy to transfer the two CLs to L2 due to the L2-L3 bandwidth of 32 B/cy; the empirical latency penalty was determined to be $T_p = 1$ cy for the 14-core HSW and $T_p = 5$ cy for the 22-core BDW. The latency penalty is strongly correlated with the number of hops in the Uncore; as BDW features more cores and each core's L3 slice forms a hop in the Uncore its latency is higher than that of the HSW chip with fewer cores/hops.

To compute the contribution of transferring the two CLs from main memory to the L3 cache, we convert the sustained memory bandwidth from GB/s to B/cy. Note that in cluster on die mode a single core can only make use of the bandwidth inside its memory domain. For the HSW, which runs at 2.3 GHz, the measured memory domain bandwidth of 32.0 GB/s corresponds to a transfer time of $64 \text{ B/CL} \cdot 2.3 \text{ GHz} / 32.0 \text{ GB/s} = 4.6 \text{ cy/CL}$ or 9.2 cy for both CLs. BDW runs at 2.1 GHz, which leads to a transfer time of $64 \text{ B/CL} \cdot 2.1 \text{ GHz} / 32.3 \text{ GB/s} = 4.2 \text{ cy/CL}$ or 8.4 cy for both CLs. The same latency penalty as for the L3 cache is applied for data coming from main memory, because the data has to be moved from the memory controller to the L3 cache segment in which the cache line is placed. The resulting ECM model inputs are $\{1 \parallel 2 \mid 2 \mid 4 + 1 \mid 9.2 + 1\}$ cy for HSW and $\{1 \parallel 2 \mid 2 \mid 4 + 5 \mid 8.4 + 5\}$ cy for BDW.

The full ECM prediction reads $\{2 \mid 4 \mid 9 \mid 19.2\}$ cy for HSW. We choose an “update” (two flops) as the basic unit of work to make performance results for different implementations comparable. The resulting unit is “updates per second” (UP/s). The expected single core performance for the HSW is thus

$$P = \frac{16 \text{ updates} \cdot 2.3 \text{ Gcy/s}}{\{2 \mid 4 \mid 9 \mid 19.2\} \text{ cy}} = \{18.40 \mid 9.20 \mid 4.09 \mid 1.92\} \text{ GUP/s} . \quad (1)$$

The predicted saturation point is at $n_S = \lceil 19.2/9.2 \rceil = 3$ cores per memory domain or 6 cores per chip. Performance at the saturation point is $P_{ECM}^S = f \cdot W_{CL} / T_{L3Mem} = 2.3 \text{ GHz} \cdot 16 \text{ updates} / 9.2 \text{ cy} = 4 \text{ GUP/s}$ per memory domain or 8 GUP/s per chip.

For BDW the full ECM prediction is $\{2 \mid 4 \mid 13 \mid 26.4\}$ cy and the resulting expected serial performance at 2.1 GHz is

$$P = \frac{16 \text{ updates} \cdot 2.1 \text{ Gcy/s}}{\{2 \mid 4 \mid 13 \mid 26.4\} \text{ cy}} = \{16.80 \mid 8.40 \mid 2.58 \mid 1.27\} \text{ GUP/s} . \quad (2)$$

The predicted saturation point is at $n_S = \lceil 26.4/8.4 \rceil = 4$ cores per memory domain or 8 cores per chip. The difference in sustained memory bandwidth between our HSW and BDW systems are marginal, so the prediction for the saturated performance is identical to that of the HSW machine.

4.1.2. Intel Xeon Phi KNC's initial many core instructions (IMCI) extensions have a SIMD width of 512 b or 64 B, corresponding to a full cache line. This means that two 512-b IMCI load instructions are needed to load the data from the L1 cache into registers, so $T_{OL} = 2$ cy. Processing the data requires one FMA instruction, which has a maximum throughput of one per cycle, resulting in $T_{OL} = 1$ cy. Note that while a KNC core is much simpler than its HSW or BDW counterpart, it is still capable of retiring two instructions in a superscalar fashion. It features two pipelines: a vector pipeline (U-pipe) with the 512-b vector processing unit attached and a scalar pipeline that handles all remaining instructions. While SIMD vector arithmetic is only possible on the U-pipe, the V-pipe

can also be used for SIMD load instructions. It is thus possible to overlap the FMA that is scheduled on the U-pipe with one of the load instructions when both instructions are paired.[†]

At a bandwidth of 32 B/cy [18], it takes $T_{L1L2} = 4$ cy to deliver the data from the L2 cache to the core. At a clock speed of 1.05 GHz and a sustained memory bandwidth of 175 GB/s the transfer time of a single CL is $64 \text{ B/CL} \cdot 1.05 \text{ GHz}/175 \text{ GB/s} = 0.4$ cy, thus 0.8 cy for both CLs. The empirically determined latency penalty for the ring interconnect amounts to $T_p = 20$ cy. The resulting ECM input is $\{1 \parallel 2 \parallel 4 \parallel 0.8 + 20\}$ cy.

The full ECM prediction is $\{2 \parallel 6 \parallel 26.8\}$ cy. It is clear from these numbers that the KNC is a strongly latency-dominated machine beyond the L2 cache. The expected performance of a single core is

$$P = \frac{16 \text{ updates} \cdot 1.05 \text{ Gcy/s}}{\{2 \parallel 6 \parallel 26.8\} \text{ cy}} = \{8.40 \parallel 2.80 \parallel 0.63\} \text{ GUP/s} . \quad (3)$$

The predicted saturation point is at $n_S = \lceil 26.8/0.8 \rceil = 34$ cores, and the maximum performance is 21.3 GUP/s.

4.1.3. IBM POWER8 In contrast to Intel Xeon and Xeon Phi chips, where cache lines are 64 B, the cache line size on the IBM PWR8 is 128 B. At a SIMD width of 16 B this means that a total of 16 VSX LOAD instructions are required to move data from the L1 cache to the registers, which takes eight cycles. The L1 cache is multi-ported, i.e., it can supply data to the registers and simultaneously receive data from the L2 cache [19]. Eight VSX FMA instructions process the data from both CLs; the kernel is thus limited by the throughput of the LOAD units and $T_{OL} = 8$ cy. As there are no non-overlapping instructions, we have $T_{nOL} = 0$ cy.

Data can be delivered from the L2 to the L1 at a bandwidth of 64 B/cy, thus $T_{L1L2} = 4$ cy. Using the documented L2-L3 bandwidth of 32 B/cy we calculate $T_{L2L3} = 8$ cy. When data is in main memory, the bandwidth of the chip-to-Centaur interconnect proves to be the bottleneck: each centaur can provide 19.2 GB/s, which translates into a peak bandwidth of 76.8 GB/s for our system. The measured sustained memory bandwidth is 73.6 GB/s, hence a CL transfer takes $128 \text{ B/CL} \cdot 2.9 \text{ GHz}/73.6 \text{ GB/s} = 5.0$ cy. Consequently, $T_{L2L4} = 10$ cy.

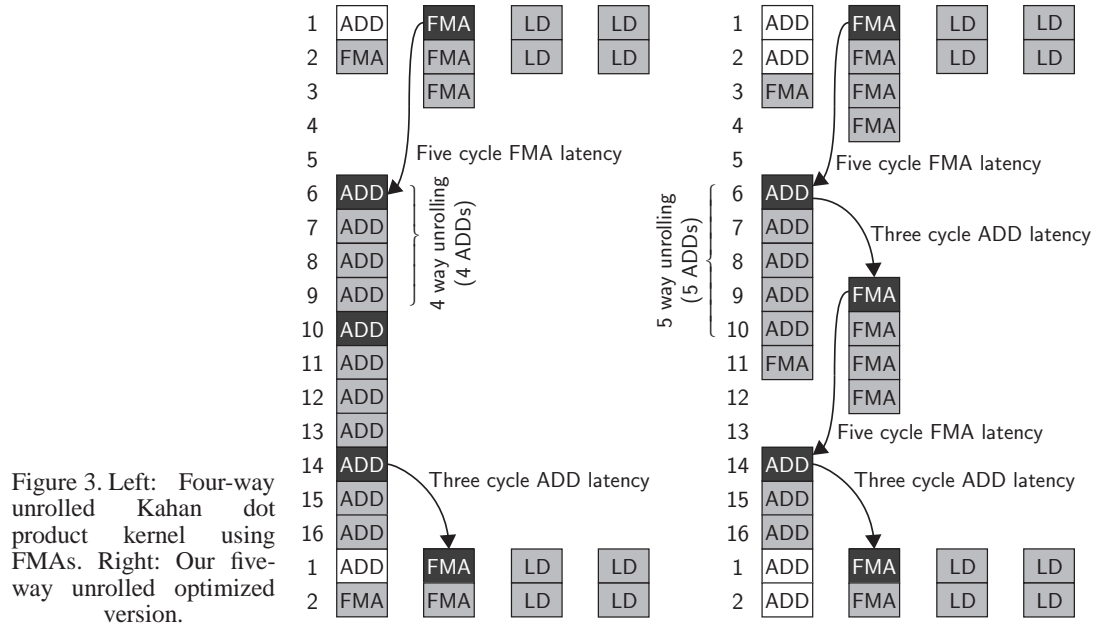
The resulting ECM input is $\{8 \parallel 0 \parallel 4 \parallel 8 \parallel 10\}$ cy. We assume a latency penalty T_p of zero, because in the measurements there is no deviation from the model prediction for data coming from the L3 cache. The reason is that on PWR8, each core has a dedicated L3 cache in which data for a particular core resides; opposed to Intel's Uncore design, no transfers across the L3 cache interconnect are necessary when accessing data from L3. The full ECM model prediction is $\{8 \parallel 8 \parallel 12 \parallel 22\}$ cy. The predicted saturation point is at $n_S = \lceil 22/10 \rceil = 3$ cores.

4.2. Kahan-enhanced scalar product

Figure 2b shows the implementation of the Kahan algorithm for the dot product. Compilers have problems with this loop code for two reasons: First, the compiler detects (correctly) a loop-carried dependency on c , which prohibits SIMD vectorization and modulo unrolling. Second, the compiler may recognize that, arithmetically, c is always equal to zero. With high optimization levels it may thus reduce the code to the naive scalar product, defeating the purpose of the Kahan algorithm. This is the reason why we use hand-coded assembly throughout this work. For comparison we also show compiler-generated Kahan code for which we ensured (by appropriate compiler options) that the algorithm is preserved.

4.2.1. Intel Haswell-EP and Broadwell-EP One iteration comprises one multiplication, four additions or subtractions, and two loads. The bottleneck on the HSW and BDW cores is thus the ADD unit (ADD and SUB are handled by the same pipeline). In the following we construct the ECM model for the AVX versions of the Kahan loop.

[†]Instruction pairing happens whenever an instruction scheduled for the U-pipe is followed directly by an instruction scheduled on the V-pipe. Restrictions about when instructions pairing can happen are complex but well documented [18].



In an AVX vectorized version of the Kahan-enhanced dot product kernel that does not use the new FMA3 extensions, we require two AVX multiplication instructions and eight AVX additions/subtractions to process one unit of work (eight scalar iterations). Multiplications can be executed speculatively several loop iterations ahead, because they have no data dependencies. This means that the five (HSW) respectively three (BDW) cycles of latency of the multiplication are not an issue. With at least four-way unrolling the add latency of three cycles can be hidden. The throughput is thus limited by the ADD unit, on which both AVX additions and subtractions are executed, resulting in $T_{OL} = 8$ cy. Because data movement is exactly the same as in the naive dot product, the remaining model inputs stay the same. This results in the following inputs for the ECM model: $\{8 \parallel 2 \parallel 2 \parallel 4 + 1 \parallel 9.2 + 1\}$ cy for HSW and $\{8 \parallel 2 \parallel 2 \parallel 4 + 5 \parallel 8.8 + 5\}$ cy for BDW. The resulting ECM predictions are $\{8 \parallel 8 \parallel 9 \parallel 19.2\}$ cy and $\{8 \parallel 8 \parallel 13 \parallel 26.8\}$ cy, respectively.

At first glance, when making use of the new FMA instructions we expect the number of in-core cycles to drop, because each core can execute two advanced vector extensions (AVX) FMA instructions per cycle. The multiplication in line four and the subtraction in line five of the source code in Fig. 2b can be handled by a single `vfmsub231ps` instruction. This reduces the number of additions/subtractions to six per cache line update so we expect T_{OL} to drop to six cycles. However, the situation is more complicated. Since the FMA instructions have `y` as input, the instruction can no longer be executed speculatively, which means that the ADD instructions now have to wait for the FMA instruction, which has a five-cycle latency on both HSW and BDW. Unfortunately, 16 addressable AVX registers are not enough to perform sufficient unrolling to completely hide this latency. It turns out that a four-way unrolled loop results in the same T_{OL} of eight cycles (see left part of Fig. 3). For a four-way unrolled kernel, intra-loop latencies play a significant role: After the first FMA has been scheduled in the first cycle of the loop (shown in black), it takes five cycles until the addition using the result of the FMA can be issued in cycle six (shown in black). The three-cycle latency of the addition is hidden by using four-way unrolling; thus it takes four cycles until the next addition corresponding to the partial sum can be retired. Finally, in the fourteenth cycle, the last addition of the first partial sum is issued. Only after the ADD latency of three cycles, the first FMA of the next loop iteration can be issued, because it uses the result of the addition as input.

Even by reusing registers that can be overwritten because their content is no longer needed, the maximum unrolling factor that can be achieved is five. Unrolling the loop with the previous strategy

```

1  vfmsub231ps zmm2, zmm0, [rdx+rax*8] # y=A[i]*B[i]-c
2  vprefetch0 [576+rsi+rax*8]          # prefetch A into L1
3
4  vaddps zmm4, zmm3, zmm2             # t=sum+y
5  vmovaps zmm0, [rsi+rax*8+64]        # load A for next iter
6
7  vsubps zmm5, zmm4, zmm3             # tmp=t-sum
8  vprefetch0 [512+rdx+rax*8]         # prefetch B into L1
9
10 vsubps zmm2, zmm5, zmm2              # c=tmp-y
11 vmovaps zmm3, zmm4                  # sum=t

```

Figure 4. Assembly code of the loop body for L2-optimized Kahan-enhanced dot product on KNC.

will result in an execution time of 18 cycles[‡] for one loop iteration (handling 2.5 cache lines at 5-way unrolling) corresponding to 7.2 cy/CL. It is possible to further decrease the runtime by “abusing” FMA operations: by keeping a vector register that has all its components set to floating-point one, we can model an addition, i.e., $y = a \times 1.0 + b$. By this optimization we can keep the loop iteration time at 16 cycles for 5-way unrolling, corresponding to a $T_{OL} = 6.4$ cy. The instruction scheduling for this version is shown on the right in Fig. 3. We replace the second unrolled additions by an FMA to increase throughput while minimizing the five-cycle latency via unrolling. The ECM model input for this optimized kernel is $\{6.4 \parallel 2 \mid 2 \mid 4 + 1 \mid 9.2 + 1\}$ cy for HSW and $\{6.4 \parallel 2 \mid 2 \mid 4 + 5 \mid 8.8 + 5\}$ cy for BDW. The resulting ECM predictions are $\{6.4 \mid 6.4 \mid 9 \mid 19.2\}$ cy and $\{6.4 \mid 6.4 \mid 13 \mid 26.8\}$ cy, respectively.

The conclusion from this analysis is that there is no expected performance difference for in-memory working sets between the naive scalar product and the Kahan version if AVX vectorization is applied to Kahan. It comes for free even in the L3 cache. Only for in-L1 and in-L2 data we expect a $2\times$ slowdown for Kahan versus the naive version even with the best possible code.

4.2.2. Intel Xeon Phi On KNC, the vector instructions performing arithmetic operations can only retire on the vector U-pipe. Thus it makes no sense to use a similar strategy as on HSW and BDW to replace additions/subtractions by fused multiply-add instructions. To process one work unit (16 scalar iterations) using 512-b SIMD instructions, the core has to execute one fused multiply-add and three additions/subtractions, yielding $T_{OL} = 4$ cy; the two 512-b loads can be executed in parallel with some of the arithmetic instructions when instructions are paired correctly, resulting in $T_{nOL} = 2$ cy. At 32 B/cy, $T_{L1L2} = 4$ cy for two cache lines. As previously determined, the sustained memory bandwidth of 175 GB/s corresponds to a transfer time of 0.4 cy/CL; thus $T_{L2Mem} = 0.8$ cy.

We found that it is necessary to use separate, specifically designed kernels to obtain the best performance for each individual cache level. The L1-optimized kernel needs no prefetching instructions at all. For data in the L2 cache, two software prefetching instructions are used, fetching eight cache lines ahead. These two instructions can be paired with arithmetic instructions and thus do not change in-core execution time (see lines two and eight in Fig. 4). For data coming from main memory, we prefetch 64 iterations ahead into the L2 cache and also keep the previous prefetching strategy of fetching cache lines eight iterations ahead from L2 into L1. The two new prefetch instructions can no longer be paired, because we run out of unpaired arithmetic instructions: The first FMA and the first ADD is paired with the LOADs that bring data into the registers; the second and third ADD/SUB are paired with the L2-L1 software prefetch instructions. The in-core execution time is thus extended by two additional cycles for the two prefetch instructions from main memory into L2. The ECM input for KNC thus is $\{4 \parallel 2 + 2_{L2} + 2_{MEM} \mid 4 \mid 0.8 + 17\}$ cy. Note that the composition of T_{nOL} is dependent on where input data is coming from: in the L1 kernel, we are retiring just two load instructions so $T_{nOL}=2$ cy; in the kernel optimized for data coming from the

[‡]By increasing the unrolling factor from four to five, we have to wait 2×5 cycles after the ADD instructions instead of 2×4 when using four-way unrolling.

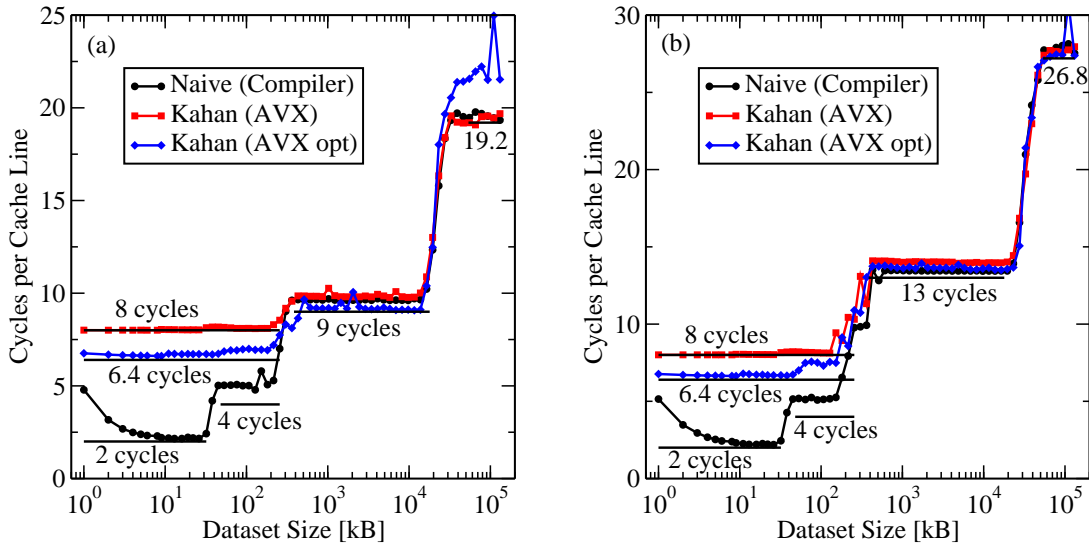


Figure 5. Single-core cycles per CL vs. data set size for AVX, AVX/FMA and the naive scalar product in SP on (a) HSW and (b) BDW. The horizontal black lines represent the ECM model predictions.

L2 cache, we need to include two prefetching instructions so $T_{nOL}=2\text{ cy}+2\text{ cy}=4\text{ cy}$; finally, for the memory-optimized kernel, we have to include two more prefetching instructions, so $T_{nOL}=6\text{ cy}$. The full ECM prediction is $\{4\mid 8\mid 27.8\}\text{ cy}$.

4.2.3. IBM POWER8 On PWR8, 16 VSX LOADs (eight 16-byte LOADs per 128-byte cache line) and required to load and an additional 32 (eight VSX FMA and 24 VSX ADD/SUB) instructions are required to process one cache line. Core throughput is limited by the two arithmetic VSX units, which require 16 cycles to process all 32 FMA/ADD/SUB instructions, resulting in $T_{OL} = 16\text{ cy}$; T_{nOL} is zero. The remaining ECM inputs are identical to the naive dot product, yielding $\{16\mid 0\mid 4\mid 8\mid 10\}$ as ECM input. The full ECM prediction is $\{16\mid 16\mid 16\mid 22\}\text{ cy}$.

5. PERFORMANCE RESULTS AND MODEL VALIDATION

5.1. Intel Haswell-EP and Broadwell-EP

Single-core benchmarking results for single precision on HSW and BDW are shown in Figs. 5a and 5b. The model describes the overall behavior very well. The naive (plain sdot) and the AVX Kahan version show identical performance in L3 cache and beyond. As predicted there is no performance drop for the AVX Kahan version from L1 to L2. The naive version as well as the AVX/FMA variant of Kahan fall short of the L2 model prediction; whether this is due to inefficiencies of the hardware prefetcher or issues with the new 64-B wide bus between L2 and L1 can only be speculated upon. We have no explanation for why the AVX/FMA optimized version shows worse in-memory performance on HSW.

In-memory scaling results on the chip level for HSW and BDW are shown in Figs. 8a and 8b. Note here that due to the cluster on die mode, the actual number of cores per memory domain is half of what the x axis in the graphs shows, i.e., the two-core run was done with one core per memory domain. This ensures that we can report the capabilities of the full chip. The number of cores required to reach saturation is underestimated in both cases. It is a well-known deficiency of the ECM model that the scaling behavior near the saturation point is not tracked correctly. We attribute this to the documented change in the prefetching strategy near memory bandwidth saturation [20]. The compiler-generated Kahan code is so slow that it misses the target of memory

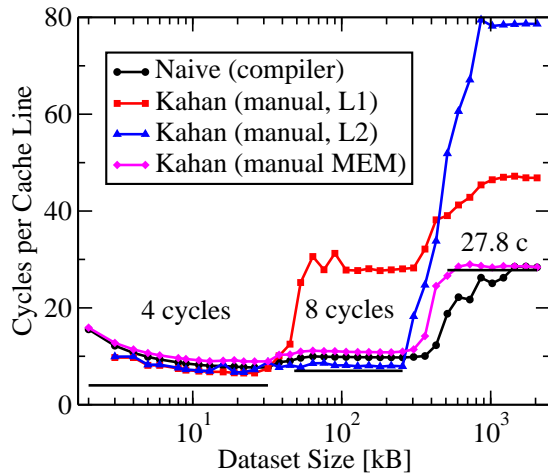


Figure 6. Single-core cycles per CL vs. data set size with different implementations tuned for specific memory hierarchy levels of the Kahan scalar product and the compiler generated naive scalar product in SP on Intel KNC. The black horizontal lines represent the ECM model predictions. All versions use 2-SMT except the manual memory-optimized kernel, which uses 4-SMT.

bandwidth saturation by far on both architectures. On HSW one would need more than twice the number of available cores to reach saturation.

5.2. Intel Xeon Phi

Figure 6 shows the SP single core results for Xeon Phi. The model fits very well as long as the special code variant in every memory hierarchy level is used. Although the Xeon Phi has a hardware prefetcher, best performance can only be achieved by appropriate software prefetching.

In-memory scaling results are shown in Fig. 8c. In accordance with Intel's guidelines, which recommend using a single thread per core when trying to reach the maximum sustained bandwidth on KNC [14], all in-memory scaling measurements were performed with 1-SMT. The compiler-generated naive and manual Kahan variants are all but identical. Xeon Phi exposes a piecewise linear scaling behavior which is not captured by the linear scaling assumption of the ECM model: Three phases can be identified, with a clear change in slope at about 20 and 50 cores. While the naive and manual Kahan codes achieve bandwidth saturation, the naive compiler version misses it by far.

5.3. IBM POWER8

Fig 7b shows the SP single core results for the PWR8 processor. The model correctly predicts the observed identical performance in L1 and L2 for the naive variant and in all memory hierarchy levels for the Kahan variant. In contrast to the Intel architectures we failed to reach the predicted instruction throughput of the processor by 20–30%. PWR8 is also more sensitive to small loop lengths. The 8 MB L3 cache is only effective up to 2 MB. Beyond this point performance dramatically decreases and fluctuates. The aggregated L4 buffer cache is not visible in the measurements. For in-memory data sets the performance improves and stabilizes. There is no documented hardware feature that could explain the erratic behavior between 2 MB and 64 MB working set size.

Fig 7a shows the impact of different SMT options on the naive sdot performance. There is no SMT setting that shows competitive performance in all memory hierarchy levels. In L1, more SMT threads lead to shorter loops and a corresponding breakdown in performance. In L2, any number of threads greater than one enables “wirespeed.” In L3 (up to 2 MB) there is clearly a strong latency effect, which can be compensated only by SMT-8. From 2 MB to the L4 capacity limit all variants exhibit the same fluctuating performance pattern with SMT-4 and SMT-8 showing the best performance. Then in memory surprisingly SMT-4 is significantly better than SMT-8. For in-memory data sets we provide two ECM model predictions: 18 cy if we assume that evicts of cachelines from L2 to L3 fully overlap with reloads from memory to L2, and 22 cy if we assume there is no overlap among those contributions. Only SMT-4 is faster than 22 cy, indicating that

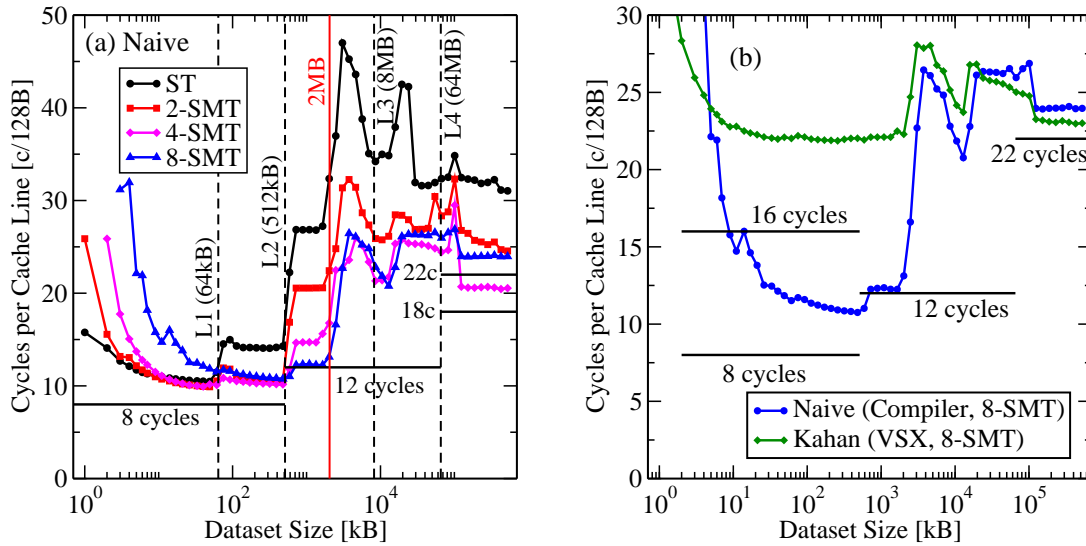


Figure 7. Single-core cycles per CL vs. data set size on PWR8. (a) Results for different SMT settings for naive scalar product using SP. (b) Comparison of compiler-generated naive scalar product and manual SIMD Kahan enhanced scalar product using SMT-8. The horizontal lines are ECM model predictions.

there is at least some overlap. More investigations are necessary to fully understand this complex behavior.

In-memory scaling results are shown in Fig. 8d. The Naive and Kahan variants show almost identical scaling behavior and quickly saturate the memory bandwidth. In contrast to the Intel architectures the compiler version of Kahan (using SMT-8) almost saturates the bandwidth.

5.4. DP performance for compiler-generated Kahan variant

As most applications rely on compiler-generated code, we show the saturation behavior of the compiler-generated Kahan variant for DP in Fig. 9. Since all compilers fail at SIMD vectorization, it is interesting to see on which architectures memory bandwidth is still achieved. On PWR8 we have already observed near-saturation in the SP case; with DP this happens at five cores. Comparing HSW and BDW, the additional cores help BDW to just about saturate whereas HSW misses this goal. KNC, as expected, misses saturation by a long shot but still achieves an absolute performance slightly better than PWR8.

5.5. Comparison across architectures

For meaningful cross-architectural comparison of the Kahan-enhanced dot product performance we report the cycles *per update* in all memory hierarchy levels (Fig. 10a) and the absolute performance for the in-memory case in GUP/s for single core as well as the full chip (Fig. 10b). In L1 and L2 all Intel architectures run close to their design specifications. PWR8 in contrast is slightly less efficient missing its design instruction throughput by 30%. In L3 and memory the results are reversed, here the Intel architectures show a significant drop in performance for L3 and also memory, especially BDW with its complex Uncore design and large number of cores, whereas PWR8 due to its lock-free memory hierarchy shows less severe performance breakdowns with increasing working set size. (Note, however, the large performance variations in a data set size window between 2 MB to over 64 MB as described in Sect. 5.3.)

Regarding absolute single-core and full-chip in-memory performance (Fig. 10b), PWR8 due to its cache architecture and higher frequency shows the best performance of all multicore chips, only surpassed by the full-chip KNC by more than a factor of two due to the latter's superior memory bandwidth.

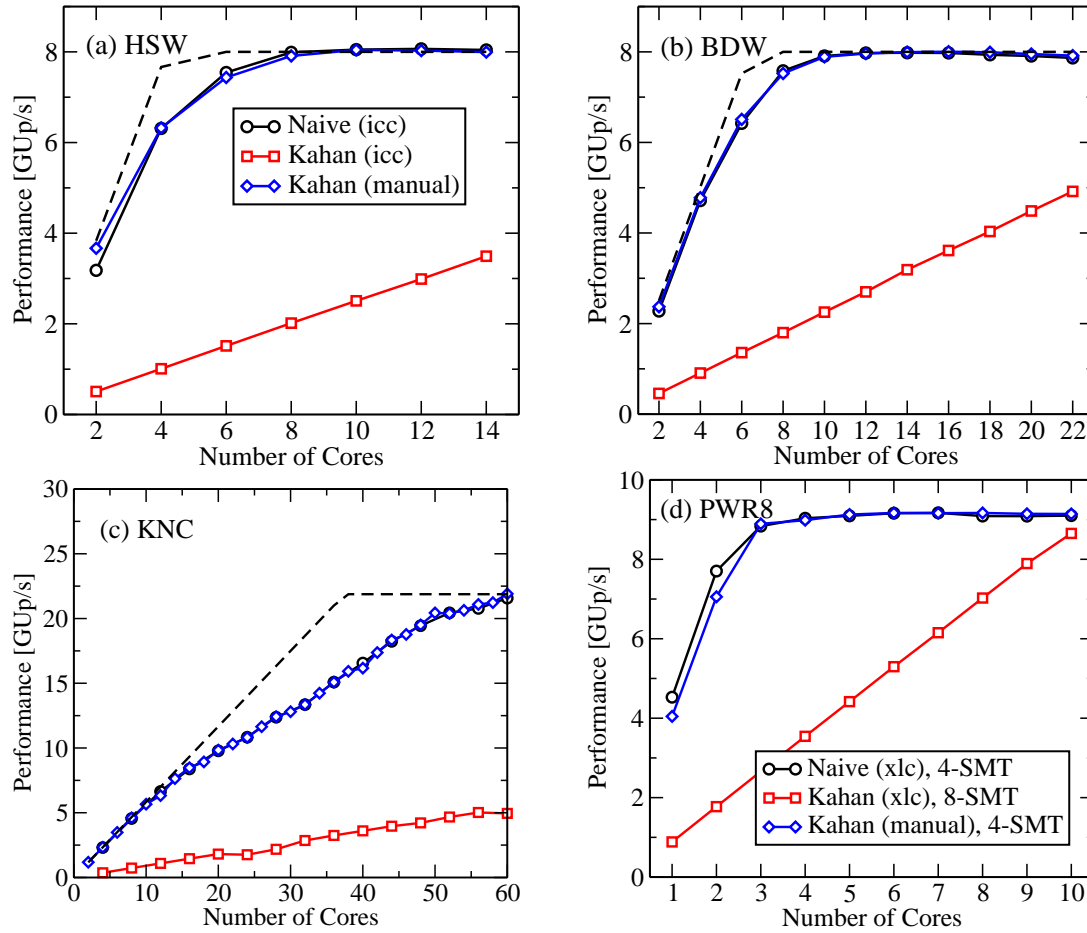


Figure 8. In-memory scaling (10 GB working set size) for different implementations of the Kahan scalar product using SP on (a) HSW, (b) BDW, (c) KNC, and (d) PWR8. One update (UP) is equivalent to five flops (one MULT, four ADDs).

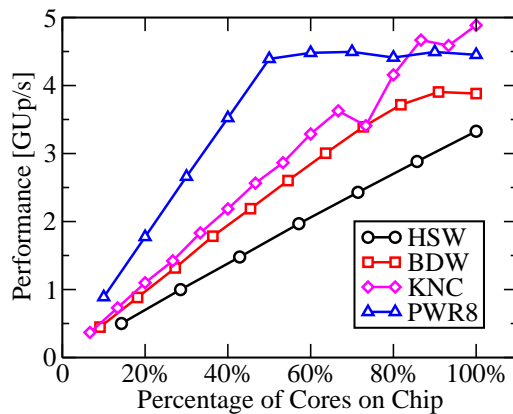


Figure 9. On-chip performance scaling of the compiler-generated Kahan-enhanced ddot on all tested processors. The saturated performance is 4 GUP/s for HSW/BDW, 10.6 GUP/s for KNC, and 4.5 GUP/s for PWR8.

6. CONCLUSION

We have investigated the performance of naive and Kahan-enhanced variants of the scalar product on a range of recent multi- and manycore chips. Using the ECM model the single-core performance in all memory hierarchy levels and the multi-core scaling for in-memory data were

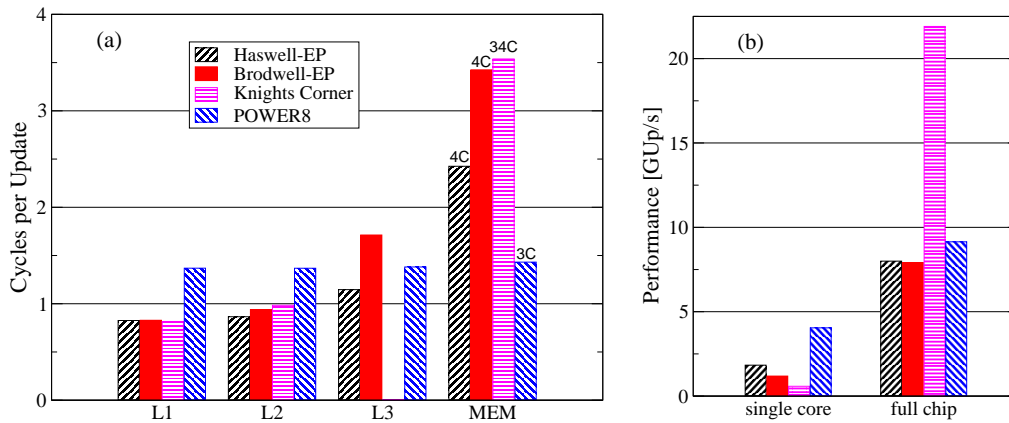


Figure 10. Comparison between all tested architectures using the manually implemented SIMD variants of the Kahan-enhanced scalar product in SP: (a) Measured single-core runtime in cycles per update in different memory hierarchy levels. The saturation point n_s is indicated above the bars for the memory-bound case (smaller is better). (b) Measured full chip performance for the in memory data set (bigger is better).

accurately described. The most important result is that even the single-threaded optimized Kahan implementation comes with no performance penalty on the Intel multicore chips under investigation compared to a naive `sdot` implementation in the L3 cache and in memory. On IBM POWER8 this applies only for in-memory data sets. On the other hand, the POWER8 is able to saturate the memory bandwidth with very few cores and provides the best single-core and chip-level performance for in-memory data. Depending on the particular architecture and whether single or double precision is used, even compiler-generated code may achieve memory bandwidth saturation on the full chip. Intel Xeon Phi as well as IBM POWER8 require special code or SMT settings to achieve best performance in different memory hierarchy levels. Further investigations are necessary to explain erratic performance behavior on POWER8 for data sets between 2 MB and 64 MB.

We emphasize that the approach and insights described here for the special case of the Kahan scalar product can serve as a blueprint for other load-dominated streaming kernels. Especially on POWER8, the ECM model still needs to be validated and adjusted using more complex codes such as stencil algorithms.

6.0.1. Acknowledgement We thank pro com and IBM Germany for access to an IBM POWER8 test system, and Intel Germany for providing an early access Broadwell-EP test system.

REFERENCES

- Goldberg D. What every computer scientist should know about floating-point arithmetic. *ACM Comput. Surv.* Mar 1991; **23**(1):5–48, DOI:10.1145/103162.103163.
- Linz P. Accurate floating-point summation. *Commun. ACM* Jun 1970; **13**(6):361–362, DOI:10.1145/362384.362498.
- Gregory J. A comparison of floating point summation methods. *Commun. ACM* Sep 1972; **15**(9), DOI:10.1145/361573.361584.
- Kahan W. Pracniques: Further remarks on reducing truncation errors. *Commun. ACM* Jan 1965; **8**(1), DOI:10.1145/363707.363723.
- Rump SM, Ogita T, Oishi S. Accurate floating-point summation part I: Faithful rounding. *SIAM J. Sci. Comput.* Oct 2008; **31**(1):189–224, DOI:10.1137/050645671.
- Zhu YK, Hayes WB. Algorithm 908: Online exact summation of floating-point streams. *ACM Trans. Math. Softw.* 2010; **37**(3):1–13, DOI:10.1145/1824801.1824815.
- Demmel J, Nguyen HD. Fast reproducible floating-point summation. *21st IEEE Symposium on Computer Arithmetic*, 2013; 163–172, DOI:10.1109/ARITH.2013.9.
- Dalton B, Wang A, Blainey B. SIMDizing pairwise sums: A summation algorithm balancing accuracy with throughput. *Proceedings of the 2014 Workshop on Programming Models for SIMD/Vector Processing, WPMVP '14*, ACM: New York, NY, USA, 2014; 65–70, DOI:10.1145/2568058.2568070.
- Hofmann J, Fey D, Riedmann M, Eitzinger J, Hager G, Wellein G. Performance analysis of the Kahan-enhanced scalar product on current multicore processors. *CoRR abs/1505.02586* 2015; URL

- <http://arxiv.org/abs/1505.02586>, accepted for PPAM'2015, the 11th International Conference on Parallel Processing and Applied Mathematics, September 6-9, 2015, Krakow, Poland.
10. Treibig J, Hager G. Introducing a performance model for bandwidth-limited loop kernels. *Parallel Processing and Applied Mathematics, Lecture Notes in Computer Science*, vol. 6067, Wyrzykowski R, Dongarra J, Karczewski K, Wasniewski J (eds.), Springer Berlin / Heidelberg, 2010; 615–624.
 11. Hager G, Treibig J, Habich J, Wellein G. Exploring performance and power properties of modern multicore chips via simple machine models. *Concurrency Computat.: Pract. Exper.* 2013; DOI: 10.1002/cpe.3180.
 12. Stengel H, Treibig J, Hager G, Wellein G. Quantifying performance bottlenecks of stencil computations using the Execution-Cache-Memory model. *Proceedings of the 29th ACM International Conference on Supercomputing, ICS '15*, ACM: New York, NY, USA, 2015, DOI:10.1145/2751205.2751240. URL <http://doi.acm.org/10.1145/2751205.2751240>.
 13. Williams S, Waterman A, Patterson D. Roofline: An insightful visual performance model for multicore architectures. *Commun. ACM* 2009; **52**(4):65–76, DOI:10.1145/1498765.1498785.
 14. Intel Corporation. Optimizing Memory Bandwidth on Stream Triad. <https://software.intel.com/en-us/articles/optimizing-memory-bandwidth-on-stream-triad>, accessed March 29, 2016.
 15. Starke WJ, Stuecheli J, Daly D, Dodson JS, Auernhammer F, Sagmeister P, Guthrie GL, Marino CF, Siegel MS, Blaner B. The cache and memory subsystems of the IBM POWER8 processor. *IBM Journal of Research and Development* 2015; **59**(1), DOI:10.1147/JRD.2014.2376131.
 16. Hofmann J, Fey D, Eitzinger J, Hager G, Wellein G. Analysis of Intel's Haswell Microarchitecture Using The ECM Model and Microbenchmarks. *CoRR abs/1511.03639* 2015; URL <http://arxiv.org/abs/1511.03639>, accepted for ARCS'2016, the 29th International Conference, April 4-7, 2016, Nuremberg, Germany.
 17. Treibig J, Hager G, Wellein G. likwid-bench: An extensible microbenchmarking platform for x86 multicore compute nodes. *Tools for High Performance Computing 2011*, Brunst H, et al. (eds.). Springer Berlin Heidelberg, 2012; 27–36, DOI:10.1007/978-3-642-31476-6_3. URL http://dx.doi.org/10.1007/978-3-642-31476-6_3.
 18. Intel Corporation. Intel Xeon Phi Core Micro-architecture. <https://software.intel.com/en-us/articles/intel-xeon-phi-core-micro-architecture>, accessed 29.3.2016.
 19. Sinharoy B, Norstrand JAV, Eickemeyer RJ, Le HQ, Leenstra J, Nguyen DQ, Konigsburg B, Ward K, Brown MD, Moreira JE, et al.. IBM POWER8 processor core microarchitecture. *IBM Journal of Research and Development* Jan 2015; **59**(1):2:1–2:21, DOI:10.1147/JRD.2014.2376112.
 20. Intel Corp. Intel64 and IA-32 Architectures Optimization Reference Manual. <http://www.intel.com/content/dam/doc/manual/64-ia-32-architectures-optimization-manual.pdf> 2015. Version: September 2015.